

# Bayesian treatment

Recap :-

- \* Our parameter estimation procedure involved:

$$\hat{\theta}_M \leftarrow \arg\max_{\theta} \log P(X|\theta)$$

- \* We used direct and indirect (EM) optimization

- \* We obtained point estimates of ' $\theta$ '

- \* What if we are interested in  $P(\theta|X)$  instead?

Why  $P(\theta|X)$ ?

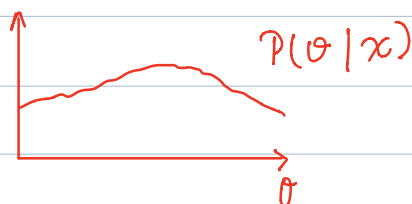
- \* Prediction using full information

$$P(z_{\text{new}} | x_{\text{new}}, X_{\text{Train}}) = \int P(z_{\text{new}} | x_{\text{new}}, \theta, X_{\text{Train}}) \cdot$$

↑  
cluster-membership

e.g. indicator for a  
new data point

$$P(\theta | x_{\text{new}}, X_{\text{Train}}) d\theta$$



\* With a point estimate, the best  
we can do is

$$P(z_{\text{new}} | x_{\text{new}}, \theta_{\text{POINT}}, X_{\text{TRAIN}})$$

\* "Integrating over  $\theta$  using  $P(\theta | x)$   
utilizes all the uncertainty/information  
about  $\theta$ "

How to compute  $P(\theta | x)$ ?

$$\text{Bayes' rule :- } P(\theta | x) = \frac{P(x | \theta) P(\theta)}{P(x)}$$

Looks simple!

But can be really hard for many  
practical problems.

$$P(\theta | x) = \frac{\overbrace{P(x | \theta)}^{\text{Likelihood}} \overbrace{P(\theta)}^{\text{Prior}}}{\underbrace{\int P(x | \theta) P(\theta) d\theta}_{P(x)}}$$

Biggest Problem

\* Normalization constant :-  $\int P(x | \theta) P(\theta) d\theta$

\* A subset of 'Likelihood, prior' pair, allows computing the normalization constant analytically.

\* For example :- Multinomial distribution

$$Pr(m | \theta, N) = \binom{N}{m_1, \dots, m_k} \prod_{k=1}^K \theta_k^{m_k}$$

Where:-

$$m := \{m_1, \dots, m_k\}$$

$$\sum \theta_k = 1, \quad \theta_k > 0 \quad \forall k$$

$$Pr(\theta | m, N) = \frac{Pr(m | \theta, N) \cdot P(\theta)}{\int Pr(m | \theta, N) \cdot P(\theta) d\theta}$$

$$\left[ \binom{N}{m_1, \dots, m_K} \prod_{k=1}^K \theta_k^{m_k} \right] P(\theta_1, \dots, \theta_K)$$

=

$$\int \binom{N}{m_1, \dots, m_K} \prod_{k=1}^K \theta_k^{m_k} P(\theta_1, \dots, \theta_K) d\theta_1 \dots d\theta_K$$

$$\theta_1, \dots, \theta_K > 0$$

$$\sum \theta_k = 1$$

\* I know of one  $P(\theta)$  which would let us compute normalization constant (with an added benefit (we'll see))

$$P(\theta_1, \dots, \theta_K | \alpha_1, \dots, \alpha_K) = \frac{\overbrace{\Gamma(\sum_k \alpha_k)}^{\text{gamma function}}}{\underbrace{\prod_k \Gamma(\alpha_k)}_{\text{hyper-parameters}}} \prod_{k=1}^K \theta_k^{\alpha_k - 1}$$

← continuous analog of factorial

Very similar to likelihood

So denominator,

using Beta Integral result

$$\binom{N}{m_1, \dots, m_K} \frac{\Gamma(\sum_k \alpha_k)}{\prod_k \Gamma(\alpha_k)} \int \prod_k \theta_k^{m_k + \alpha_k - 1} d\theta_1 \dots d\theta_K$$

$$= \binom{N}{m_1 \dots m_k} \frac{\Gamma(\sum_k d_k)}{\prod_k \Gamma(d_k)} \cdot \frac{\prod_k \Gamma(d_k + m_k)}{\Gamma(\sum_k (d_k + m_k))}$$

So

$$P(\theta | m, N, \alpha_1, \dots, \alpha_k) = \frac{\Gamma(\sum_k (d_k + m_k))}{\prod_k \Gamma(d_k + m_k)} \prod_k \theta_k^{m_k + \alpha_k - 1}$$

pseudo-counts ↓ →

{  
Same form as the  
prior

\* Hence, we can <sup>further</sup> restrict our search for prior distributions by asking one question:-

Prior  $\times$  Likelihood = Posterior



Some functional form.

\* Such priors are said to be "conjugate" to likelihood

\* In above example,  $P(\theta|\alpha)$  is a dirichlet distribution - It is conjugate to Multinomial likelihood.

\* For a family of distributions, such 'easy' priors exist - We'll generalize our treatment to that entire family (Exponential family)

Comment :- Strict adherence to conjugate priors make Bayesian treatment very easy

- However, it is a major criticism on Bayesian methods too!